

FAKE NEWS DETECTION USING SOCIAL MEDIA DATA

Industry Project Report

(TCS iON – Industry Project)

By

Mohammed Irfan Shajil .P

Campus ID: 30111

Registration Number: 23BBCACD287

Department of Computer Science

The Yenepoya Institute of Arts, Science, Commerce and Management

A Constituent Unit of Yenepoya (Deemed to be University)

Academic Year: 2023–2026

1. Abstract

The rapid growth of social media has made it easier for fake and misleading news to spread quickly. This project focuses on building a Fake News Detection system using Natural Language Processing (NLP) and Machine Learning techniques. Instead of verifying facts from external sources, the system analyzes linguistic patterns and sentiment signals present in the text. A supervised machine learning model is trained on labeled data and deployed through a web-based interface for real-time analysis.

2. Introduction

Social media platforms have become one of the primary sources of news consumption. While they enable fast information sharing, they also allow false or misleading content to reach a wide audience in a short time. Manual verification of news articles is not practical due to the volume and speed of content generation. This creates the need for automated systems that can assist in identifying fake news. This project aims to design and implement a machine learning-based solution that analyzes textual content and predicts whether a given news article is real or fake.

3. Problem Statement

The main problem addressed in this project is the identification of fake news using only textual information. The system must handle noisy and unstructured social media-style text while providing reliable predictions. The challenge lies in differentiating fake and real news without relying on external fact-checking sources.

4. Objectives

The objectives of this project are:

- To build a machine learning model for fake news classification.
- To apply NLP techniques for text preprocessing and feature extraction.
- To enhance predictions using sentiment analysis.

- To develop a simple and interactive web interface for users.
- To evaluate the system using standard performance metrics.

5. Dataset Description

The dataset used for this project consists of labeled news and social media text samples categorized as fake or real. The dataset covers multiple topics such as politics, health, and general news. Each data record contains textual content and an associated label used for supervised learning.

6. Methodology

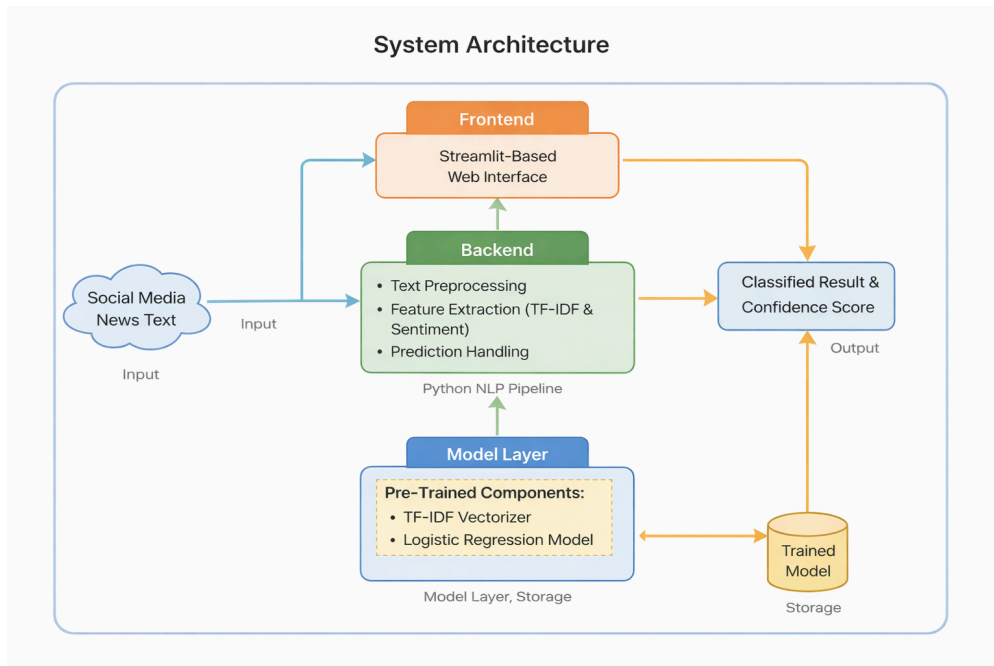
The methodology followed in this project includes the following steps:

1. Text preprocessing to remove noise and normalize input text.
2. Feature extraction using Term Frequency–Inverse Document Frequency (TF-IDF).
3. Sentiment analysis using the VADER sentiment analyzer.
4. Combining textual and sentiment features into a single feature set.
5. Training a Logistic Regression model on labeled data.
6. Applying confidence-based decision logic to handle uncertain predictions.

7. System Architecture

The system architecture consists of three main components:

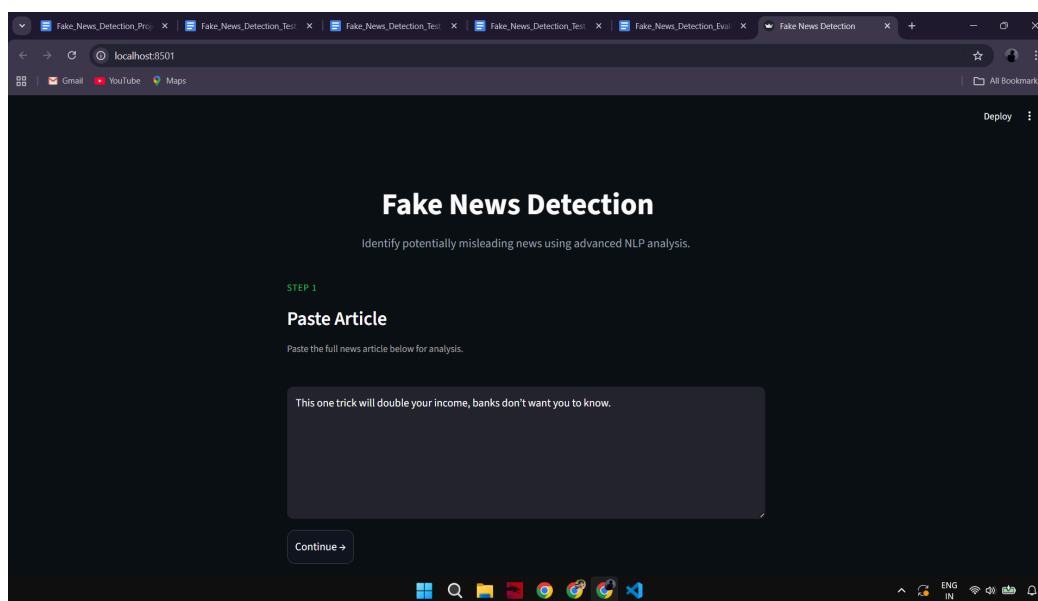
- Frontend: A Streamlit-based web interface for user input and result display.
- Backend: Python-based processing pipeline handling preprocessing, feature extraction, and prediction.
- Model Layer: Pre-trained TF-IDF vectorizer and Logistic Regression classifier loaded during runtime.



(Figure 7.1: System architecture of the Fake News Detection system)

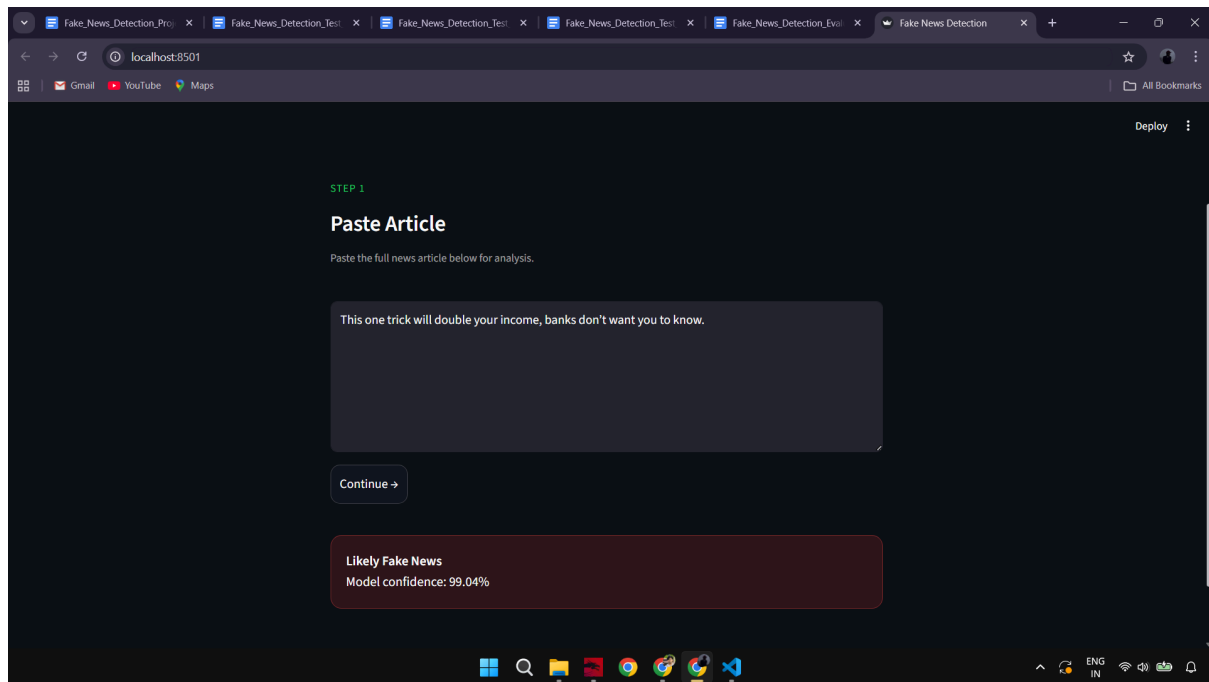
8. Implementation Details

The application is implemented using Python and libraries such as scikit-learn, NumPy, SciPy, and Streamlit. Input text is validated before processing to avoid invalid predictions. The cleaned text is vectorized using TF-IDF and combined with sentiment scores. The final prediction is generated by the trained model along with a confidence score.



(Figure 8.1: User interface for news text input and analysis)

After processing the input text, the system displays the predicted class along with a confidence score based on the trained model.



(Figure 8.2: Prediction result showing classification and confidence score)

9. Results and Evaluation

The model performance is evaluated using metrics such as accuracy, precision, recall, and F1-score. K-Fold Cross Validation is used during training to ensure better generalization. Low-confidence predictions are marked as inconclusive to improve system reliability.

10. Limitations

- The system does not perform real-time fact checking.
- Accuracy depends on the quality of the training dataset.
- Very neutral or ambiguous text may lead to inconclusive results.

11. Conclusion

This project demonstrates the effective use of NLP and machine learning techniques for fake news detection. By combining linguistic features with sentiment analysis and deploying the model through a web interface, the system provides a practical solution for identifying potentially misleading news.

12. Future Scope

Future improvements may include the use of deep learning models, integration of source credibility features, and deployment as a scalable web service with API support.